



TITLE:

Density Ratio Estimation : A Comprehensive Review (Statistical Experiment and Its Related Topics)

AUTHOR(S):

Sugiyama, Masashi; Suzuki, Taiji; Kanamori, Takafumi

CITATION:

Sugiyama, Masashi ...[et al]. Density Ratio Estimation : A Comprehensive Review (Statistical Experiment and Its Related Topics). 数理解析研究所講究録 2010, 1703: 10-31

ISSUE DATE:

2010-08

URL:

<http://hdl.handle.net/2433/170027>

RIGHT:

Density Ratio Estimation: A Comprehensive Review

Masashi Sugiyama, Tokyo Institute of Technology (sugi@cs.titech.ac.jp)

Taiji Suzuki, The University of Tokyo (s-taiji@stat.t.u-tokyo.ac.jp)

Takafumi Kanamori, Nagoya University (kanamori@is.nagoya-u.ac.jp)

Abstract

Density ratio estimation has attracted a great deal of attention in the statistics and machine learning communities since it can be used for solving various statistical data processing tasks such as non-stationarity adaptation, two-sample test, outlier detection, independence test, feature selection/extraction, independent component analysis, causal inference, and conditional probability estimation. When estimating the density ratio, it is preferable to avoid estimating densities since density estimation is known to be a hard problem. In this paper, we give a comprehensive review of density ratio estimation methods based on moment matching, probabilistic classification, and ratio matching.

1 Introduction

Recently, a new general framework of statistical data processing based on the *ratio* of probability densities has been developed (Sugiyama et al., 2009; Sugiyama et al., 2011). This density ratio framework includes various statistical data processing tasks such as non-stationarity adaptation (Shimodaira, 2000; Zadrozny, 2004; Sugiyama & Müller, 2005; Sugiyama et al., 2007; Quiñonero-Candela et al., 2009; Sugiyama et al., 2010d), outlier detection (Hido et al., 2008; Smola et al., 2009; Hido et al., 2010), change detection in time series (Kawahara & Sugiyama, 2009), conditional density estimation (Sugiyama et al., 2010c), and probabilistic classification (Sugiyama, 2010).

Furthermore, *mutual information*—which plays a central role in information theory (Cover & Thomas, 2006)—can be estimated via density ratio estimation (Suzuki et al., 2008; Suzuki et al., 2009b). Since mutual information is a measure of statistical independence between random variables, density ratio estimation can be used also for variable selection (Suzuki et al., 2009a), dimensionality reduction (Suzuki & Sugiyama, 2010), independent component analysis (Suzuki & Sugiyama, 2009), and causal inference (Yamada & Sugiyama, 2010). Thus, density ratio estimation is a promising versatile tool for statistical data processing.

A naive approach to estimating the density ratio is to separately estimate the densities corresponding to the numerator and denominator of the ratio, and then take the ratio of the estimated densities. However, this naive approach is not reliable in high-dimensional problems since division by an estimated quantity can magnify the estimation error. To overcome this drawback, various approaches to directly estimating the density ratio without going through density estimation have been explored recently, including the *moment matching approach* (Gretton et al., 2009), the *probabilistic classification approach* (Qin, 1998; Cheng & Chu, 2004; Bickel et al., 2007), and the *ratio matching approach* (Sugiyama et al., 2008; Kanamori et al., 2009a; Tsuboi et al., 2009; Yamada & Sugiyama, 2009; Yamada et al., 2010). The purpose of this paper is to provide a comprehensive review of such direct density-ratio estimation methods.

The problem of density ratio estimation addressed in this paper is formulated as follows. Let $\mathcal{X} (\subset \mathbb{R}^d)$ be the data domain, and suppose we are given independent and identically distributed (i.i.d.) samples $\{\mathbf{x}_i^{\text{nu}}\}_{i=1}^{n_{\text{nu}}}$ from a distribution with density $p_{\text{nu}}^*(\mathbf{x})$ and i.i.d. samples $\{\mathbf{x}_j^{\text{de}}\}_{j=1}^{n_{\text{de}}}$ from another distribution with density $p_{\text{de}}^*(\mathbf{x})$.

$$\{\mathbf{x}_i^{\text{nu}}\}_{i=1}^{n_{\text{nu}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{nu}}^*(\mathbf{x}) \quad \text{and} \quad \{\mathbf{x}_j^{\text{de}}\}_{j=1}^{n_{\text{de}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{de}}^*(\mathbf{x}).$$

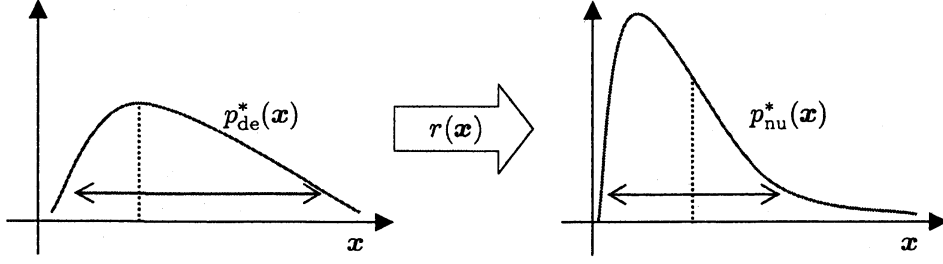


Figure 1: Matching the moments of $r(\mathbf{x})p_{\text{de}}^*(\mathbf{x})$ with those of $p_{\text{nu}}^*(\mathbf{x})$.

We assume that $p_{\text{de}}^*(\mathbf{x})$ is strictly positive over the domain \mathcal{X} . The goal is to estimate the density ratio

$$r^*(\mathbf{x}) := \frac{p_{\text{nu}}^*(\mathbf{x})}{p_{\text{de}}^*(\mathbf{x})}$$

from samples $\{\mathbf{x}_i^{\text{nu}}\}_{i=1}^{n_{\text{nu}}}$ and $\{\mathbf{x}_j^{\text{de}}\}_{j=1}^{n_{\text{de}}}$. ‘nu’ and ‘de’ indicate ‘numerator’ and ‘denominator’, respectively.

2 Moment Matching Approach

In this section, we describe the *moment matching* approach to density ratio estimation.

2.1 Preliminaries

Suppose that a one-dimensional random variable x is drawn from a probability distribution with density $p^*(x)$. Then the k -th order moment of x about the origin is defined by $\int x^k p^*(x) dx$. Note that two distributions are equivalent if and only if all moments (i.e., for $k = 1, 2, \dots$) agree with each other.

The moment matching approach to density ratio estimation tries to match the moments of $p_{\text{nu}}^*(\mathbf{x})$ and $p_{\text{de}}^*(\mathbf{x})$ via a ‘transformation’ function $r(\mathbf{x})$. More specifically, using the true density ratio $r^*(\mathbf{x})$, $p_{\text{nu}}^*(\mathbf{x})$ can be expressed as

$$p_{\text{nu}}^*(\mathbf{x}) = r^*(\mathbf{x})p_{\text{de}}^*(\mathbf{x}).$$

Thus, for a density ratio model $r(\mathbf{x})$, matching the moments of $p_{\text{nu}}^*(\mathbf{x})$ and $r(\mathbf{x})p_{\text{de}}^*(\mathbf{x})$ leads to the true density ratio $r^*(\mathbf{x})$. A schematic illustration of the moment matching approach is described in Figure 1.

2.2 Finite-Order Approach

The simplest implementation of moment matching would be to match the first-order moment (i.e., the mean):

$$\underset{r}{\operatorname{argmin}} \left\| \int \mathbf{x} r(\mathbf{x}) p_{\text{de}}^*(\mathbf{x}) d\mathbf{x} - \int \mathbf{x} p_{\text{nu}}^*(\mathbf{x}) d\mathbf{x} \right\|^2,$$

where $\|\cdot\|$ denotes the Euclidean norm. Its non-linear variant can be obtained using some non-linear function $\phi(\mathbf{x}) : \mathbb{R}^d \rightarrow \mathbb{R}$ as

$$\underset{r}{\operatorname{argmin}} \left(\int \phi(\mathbf{x}) r(\mathbf{x}) p_{\text{de}}^*(\mathbf{x}) d\mathbf{x} - \int \phi(\mathbf{x}) p_{\text{nu}}^*(\mathbf{x}) d\mathbf{x} \right)^2.$$

This non-linear method can be easily extended to multiple components by using a vector-valued function $\phi(\mathbf{x}) : \mathbb{R}^d \rightarrow \mathbb{R}^m$ as

$$\operatorname{argmin}_r \text{MM}'(r), \quad \text{where} \quad \text{MM}'(r) := \left\| \int \phi(\mathbf{x}) r(\mathbf{x}) p_{\text{de}}^*(\mathbf{x}) d\mathbf{x} - \int \phi(\mathbf{x}) p_{\text{nu}}^*(\mathbf{x}) d\mathbf{x} \right\|^2,$$

where ‘MM’ stands for ‘moment matching’. Let us ignore the irrelevant constant in $\text{MM}'(r)$, and define the rest as $\text{MM}(r)$:

$$\text{MM}(r) := \left\| \int \phi(\mathbf{x}) r(\mathbf{x}) p_{\text{de}}^*(\mathbf{x}) d\mathbf{x} \right\|^2 - 2 \left\langle \int \phi(\mathbf{x}) r(\mathbf{x}) p_{\text{de}}^*(\mathbf{x}) d\mathbf{x}, \int \phi(\mathbf{x}) p_{\text{nu}}^*(\mathbf{x}) d\mathbf{x} \right\rangle,$$

where $\langle \cdot, \cdot \rangle$ denotes the inner product.

In practice, the expectations over $p_{\text{nu}}^*(\mathbf{x})$ and $p_{\text{de}}^*(\mathbf{x})$ in $\text{MM}(r)$ are replaced by sample averages. That is, for an n_{de} -dimensional vector $\mathbf{r}_{\text{de}}^* := (r^*(\mathbf{x}_1^{\text{de}}), \dots, r^*(\mathbf{x}_{n_{\text{de}}}^{\text{de}}))^{\top}$ where \top denotes the transpose, an estimator $\hat{\mathbf{r}}_{\text{de}}$ of \mathbf{r}_{de}^* can be obtained by solving the following optimization problem.

$$\hat{\mathbf{r}}_{\text{de}} := \operatorname{argmin}_{\mathbf{r} \in \mathbb{R}^{n_{\text{de}}}} \widehat{\text{MM}}(\mathbf{r}), \quad \text{where} \quad \widehat{\text{MM}}(\mathbf{r}) := \frac{1}{n_{\text{de}}^2} \mathbf{r}^{\top} \Phi_{\text{de}}^{\top} \Phi_{\text{de}} \mathbf{r} - \frac{2}{n_{\text{de}} n_{\text{nu}}} \mathbf{r}^{\top} \Phi_{\text{de}}^{\top} \Phi_{\text{nu}} \mathbf{1}_{n_{\text{nu}}}. \quad (1)$$

$\mathbf{1}_n$ denotes the n -dimensional vector with all ones. Φ_{nu} and Φ_{de} are the $t \times n_{\text{nu}}$ and $t \times n_{\text{de}}$ design matrices defined by $\Phi_{\text{nu}} := (\phi(\mathbf{x}_1^{\text{nu}}), \dots, \phi(\mathbf{x}_{n_{\text{nu}}}^{\text{nu}}))$ and $\Phi_{\text{de}} := (\phi(\mathbf{x}_1^{\text{de}}), \dots, \phi(\mathbf{x}_{n_{\text{de}}}^{\text{de}}))$, respectively. Taking the derivative of the objective function (1) with respect to \mathbf{r} and setting it to zero, we have

$$\frac{2}{n_{\text{de}}^2} \Phi_{\text{de}}^{\top} \Phi_{\text{de}} \mathbf{r} - \frac{2}{n_{\text{de}} n_{\text{nu}}} \Phi_{\text{de}}^{\top} \Phi_{\text{nu}} \mathbf{1}_{n_{\text{nu}}} = \mathbf{0}_t,$$

where $\mathbf{0}_t$ denotes the t -dimensional vector with all zeros. Solving this equation with respect to \mathbf{r} , we can obtain the solution analytically as

$$\hat{\mathbf{r}}_{\text{de}} = \frac{n_{\text{de}}}{n_{\text{nu}}} (\Phi_{\text{de}}^{\top} \Phi_{\text{de}})^{-1} \Phi_{\text{de}}^{\top} \Phi_{\text{nu}} \mathbf{1}_{n_{\text{nu}}}.$$

One may add a normalization constraint $\frac{1}{n_{\text{de}}} \mathbf{1}_{n_{\text{de}}}^{\top} \mathbf{r} = 1$ to the optimization problem (1). Then the optimization problem becomes a linearly constrained quadratic program. Thus, a numerical solver may be needed to compute the solution. Furthermore, a non-negativity constraint $\mathbf{r} \geq \mathbf{0}_{n_{\text{de}}}$ and/or an upper bound for a positive constant B (i.e., $\mathbf{r} \leq B \mathbf{1}_{n_{\text{de}}}$) may also be incorporated in the optimization problem (1), where inequalities for vectors are applied in the element-wise manner. Even with these modifications, the optimization problem is still a linearly constrained quadratic program, so its solution can be numerically computed by standard optimization software.

The above moment-matching method gives an estimate of the density ratio values at the denominator sample points $\{\mathbf{x}_j^{\text{de}}\}_{j=1}^{n_{\text{de}}}$. If one wants to estimate the entire ratio function $r^*(\mathbf{x})$, the following linear density-ratio model may be used instead (Kanamori et al., 2009b):

$$r(\mathbf{x}) = \psi(\mathbf{x})^{\top} \boldsymbol{\theta}, \quad (2)$$

where $\psi(\mathbf{x}) : \mathbb{R}^d \rightarrow \mathbb{R}^b$ is a basis function vector and $\boldsymbol{\theta} (\in \mathbb{R}^b)$ is a parameter vector. We assume that the basis functions are non-negative: $\psi(\mathbf{x}) \geq \mathbf{0}_b$. Then model outputs at $\{\mathbf{x}_j^{\text{de}}\}_{j=1}^{n_{\text{de}}}$ are expressed in terms of the parameter vector $\boldsymbol{\theta}$ as

$$(r(\mathbf{x}_1^{\text{de}}), \dots, r(\mathbf{x}_{n_{\text{de}}}^{\text{de}}))^{\top} = \Psi_{\text{de}}^{\top} \boldsymbol{\theta},$$

where Ψ_{de} is the $b \times n_{\text{de}}$ design matrix defined by $\Psi_{\text{de}} := (\psi(\mathbf{x}_1^{\text{de}}), \dots, \psi(\mathbf{x}_{n_{\text{de}}}^{\text{de}}))$. Then, following Eq.(1), the parameter θ is learned as follows.

$$\hat{\theta} := \underset{\theta \in \mathbb{R}^b}{\operatorname{argmin}} \left[\frac{1}{n_{\text{de}}^2} \theta^\top \Psi_{\text{de}} \Phi_{\text{de}}^\top \Phi_{\text{de}} \Psi_{\text{de}}^\top \theta - \frac{2}{n_{\text{de}} n_{\text{nu}}} \theta^\top \Psi_{\text{de}} \Phi_{\text{de}}^\top \Phi_{\text{nu}} \mathbf{1}_{n_{\text{nu}}} \right]. \quad (3)$$

Taking the derivative of the above objective function with respect to θ and setting it to zero, we have the solution $\hat{\theta}$ analytically as

$$\hat{\theta} = \frac{n_{\text{de}}}{n_{\text{nu}}} (\Psi_{\text{de}} \Phi_{\text{de}}^\top \Phi_{\text{de}} \Psi_{\text{de}}^\top)^{-1} \Psi_{\text{de}} \Phi_{\text{de}}^\top \Phi_{\text{nu}} \mathbf{1}_{n_{\text{nu}}}.$$

One may include a normalization constraint, a non-negativity constraint (given that the basis function is non-negative), and a regularization constraint to the optimization problem (3):

$$\frac{1}{n_{\text{de}}} \mathbf{1}_{n_{\text{de}}}^\top \Psi_{\text{de}}^\top \theta = 1, \quad \theta \geq \mathbf{0}_b, \quad \text{and} \quad \theta \leq B \mathbf{1}_b.$$

Then the optimization problem becomes a linearly constrained quadratic program, whose solution can be obtained by a standard numerical solver.

The upper-bound parameter B , which works as a regularizer, may be optimized by *cross-validation* (CV). That is, the numerator and denominator samples $D^{\text{nu}} = \{\mathbf{x}_i^{\text{nu}}\}_{i=1}^{n_{\text{nu}}}$ and $D^{\text{de}} = \{\mathbf{x}_j^{\text{de}}\}_{j=1}^{n_{\text{de}}}$ are first divided into T disjoint subsets $\{D_t^{\text{nu}}\}_{t=1}^T$ and $\{D_t^{\text{de}}\}_{t=1}^T$, respectively. Then a density ratio estimator $\hat{r}_t(\mathbf{x})$ is obtained from $D^{\text{nu}} \setminus D_t^{\text{nu}}$ and $D^{\text{de}} \setminus D_t^{\text{de}}$ (i.e., all samples without D_t^{nu} and D_t^{de}), and its moment matching error is computed for the hold-out samples D_t^{nu} and D_t^{de} .

$$\widetilde{\text{MM}}_t(\hat{r}) := \left(\frac{1}{|D_t^{\text{de}}|} \sum_{\mathbf{x}^{\text{de}} \in D_t^{\text{de}}} \phi(\mathbf{x}^{\text{de}}) \hat{r}_t(\mathbf{x}^{\text{de}}) \right)^2 - \frac{2}{|D_t^{\text{de}}| |D_t^{\text{nu}}|} \left(\sum_{\mathbf{x}^{\text{de}} \in D_t^{\text{de}}} \phi(\mathbf{x}^{\text{de}}) \hat{r}_t(\mathbf{x}^{\text{de}}) \right)^\top \left(\sum_{\mathbf{x}^{\text{nu}} \in D_t^{\text{nu}}} \phi(\mathbf{x}^{\text{nu}}) \right),$$

where $|D|$ denotes the number of elements in the set D . This procedure is repeated for $t = 1, \dots, T$, and the average of the above hold-out moment matching error over all t is computed.

$$\widetilde{\text{MM}} := \frac{1}{T} \sum_{t=1}^T \widetilde{\text{MM}}_t.$$

Then the upper-bound parameter B that minimizes $\widetilde{\text{MM}}$ is chosen. Availability of CV would be one of the advantages of the inductive method (i.e., learning the entire density-ratio function).

2.3 Infinite-Order Approach

Matching a finite number of moments does not necessarily result in the true density ratio function $r^*(\mathbf{x})$, even if infinitely many samples are available. In order to guarantee that the true density ratio function can always be obtained in the large-sample limit, all moments up to the infinite order need to be matched.

Kernel mean matching (KMM) allows one to efficiently match all the moments (Gretton et al., 2009). The basic idea of KMM is common to the above finite-order approach, but a *universal reproducing kernel* $K(\mathbf{x}, \mathbf{x}')$ (Steinwart, 2001) is used as a non-linear transformation. The *Gaussian kernel*

$$K(\mathbf{x}, \mathbf{x}') = \exp \left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2} \right) \quad (4)$$

is an example of universal reproducing kernels. It has been shown that the solution of the following optimization problem agrees with the true density ratio (Gretton et al., 2009):

$$\min_{r \in \mathcal{H}} \left\| \int K(\mathbf{x}, \cdot) p_{\text{nu}}^*(\mathbf{x}) d\mathbf{x} - \int K(\mathbf{x}, \cdot) r(\mathbf{x}) p_{\text{de}}^*(\mathbf{x}) d\mathbf{x} \right\|_{\mathcal{H}}^2,$$

where \mathcal{H} denotes a universal reproducing kernel Hilbert space and $\|\cdot\|_{\mathcal{H}}$ denotes its norm.

An empirical version of the above problem is reduced to

$$\min_{\mathbf{r} \in \mathbb{R}^{n_{\text{de}}}} \left[\frac{1}{n_{\text{de}}^2} \mathbf{r}^\top \mathbf{K}_{\text{de,de}} \mathbf{r} - \frac{2}{n_{\text{de}} n_{\text{nu}}} \mathbf{r}^\top \mathbf{K}_{\text{de,nu}} \mathbf{1}_{n_{\text{nu}}} \right],$$

where $\mathbf{K}_{\text{de,nu}}$ and $\mathbf{K}_{\text{de,de}}$ denote the Gram matrices defined by $[\mathbf{K}_{\text{de,nu}}]_{i,j} = K(\mathbf{x}_i^{\text{de}}, \mathbf{x}_j^{\text{nu}})$ and $[\mathbf{K}_{\text{de,de}}]_{i,i'} = K(\mathbf{x}_i^{\text{de}}, \mathbf{x}_{i'}^{\text{de}})$, respectively. In the same way as the finite-order case, the solution can be obtained analytically as

$$\hat{\mathbf{r}}_{\text{de}} = \frac{n_{\text{de}}}{n_{\text{nu}}} \mathbf{K}_{\text{de,de}}^{-1} \mathbf{K}_{\text{de,nu}} \mathbf{1}_{n_{\text{nu}}}. \quad (5)$$

If necessary, one may include a non-negativity constraint, a normalization constraint, and an upper bound in the same way as the finite-order case. Then the solution can be numerically obtained by solving a linearly constrained quadratic programming problem.

For the linear density-ratio model (2), an inductive variant of KMM is formulated as

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^b} \left[\frac{1}{n_{\text{de}}^2} \boldsymbol{\theta}^\top \boldsymbol{\Psi}_{\text{de}} \mathbf{K}_{\text{de,de}} \boldsymbol{\Psi}_{\text{de}}^\top \boldsymbol{\theta} - \frac{2}{n_{\text{de}} n_{\text{nu}}} \boldsymbol{\theta}^\top \boldsymbol{\Psi}_{\text{de}} \mathbf{K}_{\text{de,nu}} \mathbf{1}_{n_{\text{nu}}} \right].$$

As shown above, KMM utilizes universal reproducing kernels such as the Gaussian kernel (4) to efficiently match all the moments. Theoretically, KMM is consistent for any universal reproducing kernels. However, its practical performance heavily depends on the choice of kernels such as the Gaussian width σ , and such kernel parameters cannot be simply optimized by cross-validation even in the induction cases. This is because one is not finding a Gaussian width value that matches the moments well. Thus, optimizing σ over the moment matching criterion may not be appropriate as a model selection strategy. A popular heuristic is to use the *median distance* between samples as the Gaussian width σ (Schölkopf & Smola, 2002). However, there seems no strong justification for this heuristic.

2.4 Remarks

Density ratio estimation by moment matching can successfully avoid density estimation.

The finite-order moment matching method (Section 2.2) is simple and computationally efficient, if the number of matching moments is kept reasonably small. However, the finite-order approach is not necessarily consistent. On the other hand, the infinite-order moment matching method (Section 2.3), *kernel mean matching* (KMM), can efficiently match all the moments by making use of universal reproducing kernels. Indeed, KMM has an excellent theoretical property that it is consistent. However, KMM has a limitation in model selection—there is no known method for determining the kernel parameter (such as the Gaussian kernel width). A popular heuristic of setting the Gaussian width to the median distance between samples would be useful in some cases, but this is perhaps not always reasonable.

In this section, moment matching was performed in terms of the squared norm, which led to an analytic-form solution (if no constraint is imposed). As shown in Kanamori et al. (2009b),

moment matching can be generalized to various divergences. Such a generalized KMM method actually has a close connection with the *ratio matching approach* explained in Section 4. However, the ratio matching approach is more preferable than the generalized KMM approach because of the following reasons:

- The ratio matching approach is equipped with a natural CV procedure for model selection (Sugiyama et al., 2008; Kanamori et al., 2009a). Thus no heuristic is required for choosing the Gaussian width.
- The ratio matching approach is proved to be numerically more stable than the KMM approach in terms of *condition numbers* (Kanamori et al., 2009b).

3 Probabilistic Classification Approach

In this section, we describe a framework of density ratio estimation through *probabilistic classification*.

3.1 Preliminaries

The basic idea of the probabilistic classification approach is to learn a probabilistic classifier which separates samples $\{\mathbf{x}_i^{\text{nu}}\}_{i=1}^{n_{\text{nu}}}$ drawn from $p_{\text{nu}}^*(\mathbf{x})$ and samples $\{\mathbf{x}_j^{\text{de}}\}_{j=1}^{n_{\text{de}}}$ drawn from $p_{\text{de}}^*(\mathbf{x})$ (Qin, 1998; Cheng & Chu, 2004; Bickel et al., 2007).

Let us assign a label $y = +1$ to $\{\mathbf{x}_i^{\text{nu}}\}_{i=1}^{n_{\text{nu}}}$ and $y = -1$ to $\{\mathbf{x}_j^{\text{de}}\}_{j=1}^{n_{\text{de}}}$, respectively. Then the two densities are written as $p_{\text{nu}}^*(\mathbf{x}) = p^*(\mathbf{x}|y = +1)$ and $p_{\text{de}}^*(\mathbf{x}) = p^*(\mathbf{x}|y = -1)$, respectively. Note that y is regarded as a random variable here. An application of Bayes' theorem,

$$p^*(\mathbf{x}|y) = \frac{p^*(y|\mathbf{x})p^*(\mathbf{x})}{p^*(y)},$$

yields that the density ratio can be expressed in terms of y as follows:

$$\begin{aligned} r^*(\mathbf{x}) &= \frac{p_{\text{nu}}^*(\mathbf{x})}{p_{\text{de}}^*(\mathbf{x})} = \left(\frac{p^*(y = +1|\mathbf{x})p^*(\mathbf{x})}{p^*(y = +1)} \right) \left(\frac{p^*(y = -1|\mathbf{x})p^*(\mathbf{x})}{p^*(y = -1)} \right)^{-1} \\ &= \frac{p^*(y = -1)p^*(y = +1|\mathbf{x})}{p^*(y = +1)p^*(y = -1|\mathbf{x})}. \end{aligned}$$

The ratio $p^*(y = -1)/p^*(y = +1)$ may be simply approximated by the ratio of the number of samples:

$$\frac{p^*(y = -1)}{p^*(y = +1)} \approx \frac{n_{\text{de}}/(n_{\text{de}} + n_{\text{nu}})}{n_{\text{nu}}/(n_{\text{de}} + n_{\text{nu}})} = \frac{n_{\text{de}}}{n_{\text{nu}}}.$$

The 'class' posterior probability $p^*(y|\mathbf{x})$ may be approximated by separating $\{\mathbf{x}_i^{\text{nu}}\}_{i=1}^{n_{\text{nu}}}$ and $\{\mathbf{x}_j^{\text{de}}\}_{j=1}^{n_{\text{de}}}$ using a probabilistic classifier. Thus, given an estimator of the class posterior probability, $\hat{p}(y|\mathbf{x})$, a density ratio estimator $\hat{r}(\mathbf{x})$ can be constructed as

$$\hat{r}(\mathbf{x}) = \frac{n_{\text{de}} \hat{p}(y = +1|\mathbf{x})}{n_{\text{nu}} \hat{p}(y = -1|\mathbf{x})}. \quad (6)$$

A practical advantage of the probabilistic classification approach would be its easy implementability. Indeed, one can directly use standard classification algorithms for density ratio estimation. Below, an example of probabilistic classifiers, logistic regression, is described. For making the explanation simple, we consider a set of paired samples, $\{(\mathbf{x}_k, y_k)\}_{k=1}^n$, where $(\mathbf{x}_1, \dots, \mathbf{x}_n) := (\mathbf{x}_1^{\text{nu}}, \dots, \mathbf{x}_{n_{\text{nu}}}^{\text{nu}}, \mathbf{x}_1^{\text{de}}, \dots, \mathbf{x}_{n_{\text{de}}}^{\text{de}})$ and $(y_1, \dots, y_n) := (+1, \dots, +1, -1, \dots, -1)$ for $n = n_{\text{nu}} + n_{\text{de}}$.

3.2 Logistic Regression Classifier

Here a popular classification algorithm called *logistic regression* (LR) (Hastie et al., 2001) is explained.

The LR classifier employs a parametric model of the following form for expressing the class-posterior probability $p^*(y|\mathbf{x})$,

$$p(y|\mathbf{x}; \boldsymbol{\theta}) = \left(1 + \exp\left(-y\boldsymbol{\psi}(\mathbf{x})^\top \boldsymbol{\theta}\right)\right)^{-1},$$

where $\boldsymbol{\psi}(\mathbf{x}) : \mathbb{R}^d \rightarrow \mathbb{R}^b$ is a basis function vector and $\boldsymbol{\theta} (\in \mathbb{R}^b)$ is a parameter vector. The parameter vector $\boldsymbol{\theta}$ is learned so that the penalized log-likelihood is maximized, which is equivalently expressed as follows:

$$\hat{\boldsymbol{\theta}} := \underset{\boldsymbol{\theta} \in \mathbb{R}^m}{\operatorname{argmin}} \left[\sum_{k=1}^n \log \left(1 + \exp\left(-y_k \boldsymbol{\psi}(\mathbf{x}_k)^\top \boldsymbol{\theta}\right)\right) + \lambda \boldsymbol{\theta}^\top \boldsymbol{\theta} \right], \quad (7)$$

where $\lambda \boldsymbol{\theta}^\top \boldsymbol{\theta}$ is a penalty term included for regularization purposes.

Since the objective function in Eq.(7) is convex, the global optimal solution can be obtained by a standard non-linear optimization technique such as the *gradient descent method* or (*quasi-Newton methods*) (Hastie et al., 2001; Minka, 2007). An LR model classifies a new input sample \mathbf{x} by choosing the most probable class:

$$\hat{y} = \underset{y=\pm 1}{\operatorname{argmax}} p(y|\mathbf{x}; \hat{\boldsymbol{\theta}}). \quad (8)$$

Finally, a density ratio estimator $\hat{r}_{\text{LR}}(\mathbf{x})$ is given by

$$\begin{aligned} \hat{r}_{\text{LR}}(\mathbf{x}) &= \frac{n_{\text{de}}}{n_{\text{nu}}} \frac{1 + \exp\left(\boldsymbol{\psi}(\mathbf{x})^\top \hat{\boldsymbol{\theta}}\right)}{1 + \exp\left(-\boldsymbol{\psi}(\mathbf{x})^\top \hat{\boldsymbol{\theta}}\right)} = \frac{n_{\text{de}}}{n_{\text{nu}}} \frac{\exp\left(\boldsymbol{\psi}(\mathbf{x})^\top \hat{\boldsymbol{\theta}}\right) \left\{ \exp\left(-\boldsymbol{\psi}(\mathbf{x})^\top \hat{\boldsymbol{\theta}}\right) + 1 \right\}}{1 + \exp\left(-\boldsymbol{\psi}(\mathbf{x})^\top \hat{\boldsymbol{\theta}}\right)} \\ &= \frac{n_{\text{de}}}{n_{\text{nu}}} \exp\left(\boldsymbol{\psi}(\mathbf{x})^\top \hat{\boldsymbol{\theta}}\right). \end{aligned}$$

When multi-class LR classifiers are used, density ratios among multiple densities can be estimated simultaneously. This is useful, e.g., for solving *multi-task learning* problems (Bickel et al., 2008).

When the LR model is *correctly specified*, i.e., there exists $\boldsymbol{\theta}^*$ such that $p(y|\mathbf{x}; \boldsymbol{\theta}^*) = p^*(y|\mathbf{x})$, the LR approach is optimal among a class of semi-parametric estimators in the sense that the asymptotic variance is minimized (Qin, 1998). However, when the model is misspecified (which would be the case in practice), the ratio-matching approach explained in Section 4 is more preferable (Kanamori et al., 2010).

3.3 Model Selection by Cross-Validation

An important advantage of the probabilistic classification approach is that model selection (i.e., tuning the basis functions and the regularization parameter) is possible by standard CV, since the learning problem involved in this framework is a standard supervised classification problem.

More specifically, the numerator and denominator samples $D^{\text{nu}} = \{\mathbf{x}_i^{\text{nu}}\}_{i=1}^{n_{\text{nu}}}$ and $D^{\text{de}} = \{\mathbf{x}_j^{\text{de}}\}_{j=1}^{n_{\text{de}}}$ are divided into T disjoint subsets $\{D_t^{\text{nu}}\}_{t=1}^T$ and $\{D_t^{\text{de}}\}_{t=1}^T$, respectively. Then a probabilistic classifier $\hat{p}_t(y|\mathbf{x})$ is obtained using $D^{\text{nu}} \setminus D_t^{\text{nu}}$ and $D^{\text{de}} \setminus D_t^{\text{de}}$ (i.e., all samples without

D_t^{nu} and D_t^{de}), and its misclassification error (ME) for the hold-out samples D_t^{nu} and D_t^{de} is computed:

$$\widetilde{\text{ME}}_t := \frac{1}{|D_t^{\text{nu}}|} \sum_{\mathbf{x}^{\text{nu}} \in D_t^{\text{nu}}} I(\arg\max_{y=\pm 1} \hat{p}_t(y|\mathbf{x}^{\text{nu}}) = +1) + \frac{1}{|D_t^{\text{de}}|} \sum_{\mathbf{x}^{\text{de}} \in D_t^{\text{de}}} I(\arg\max_{y=\pm 1} \hat{p}_t(y|\mathbf{x}^{\text{de}}) = -1),$$

where $I(\cdot)$ is the indicator function: $I(c) = 1$ if c is true and $I(c) = 0$ otherwise. This procedure is repeated for $t = 1, \dots, T$, and the average misclassification error over all t is computed.

$$\widetilde{\text{ME}} := \frac{1}{T} \sum_{t=1}^T \widetilde{\text{ME}}_t.$$

Then the model that minimizes $\widetilde{\text{ME}}$ is chosen.

3.4 Remarks

Density ratio estimation by probabilistic classification can successfully avoid density estimation by casting the problem of density ratio estimation as the problem of learning the ‘class’ posterior probability. An advantage of the probabilistic classification approach over the moment matching approach explained in Section 2 is that CV can be used for model selection. Furthermore, existing software packages of classification algorithms can be directly used for density ratio estimation.

As shown in Qin (1998), the probabilistic classification approach with LR has a suitable property: if the LR model is *correctly specified*, the probabilistic classification approach is optimal among a broad class of semi-parametric estimators. However, this strong theoretical property is not true when the correct model assumption is not fulfilled. In such cases, the ratio-matching approach explained in Section 4 is more preferable (Kanamori et al., 2010).

4 Ratio Matching Approach

In this section, we describe the *ratio matching* approach to density ratio estimation.

4.1 Preliminaries

A basic idea of ratio matching is to directly match a density ratio model $r(\mathbf{x})$ to the true density ratio $r^*(\mathbf{x})$ under some divergence (Figure 2). At a glance, the ratio matching problem is equivalent to the regression problem. However, ratio matching is essentially different from regression since samples of the true ratio are not available. Here, we employ the *Bregman (BR) divergence* for measuring the discrepancy between the true density ratio and the density ratio model (Bregman, 1967).

The BR divergence is an extension of the Euclidean distance to a class of distances that all share similar properties. Let f be a differentiable and strictly convex function. Then the BR divergence associated with f from t^* to t is defined as

$$\text{BR}'_f(t^*||t) := f(t^*) - f(t) - \nabla f(t)(t^* - t).$$

Note that $f(t) + \nabla f(t)(t^* - t)$ is the value of the first-order *Taylor expansion* of f around point t evaluated at point t^* . Thus, the BR divergence evaluates the difference between the value of f at point t^* and its linear extrapolation from t (Figure 3).

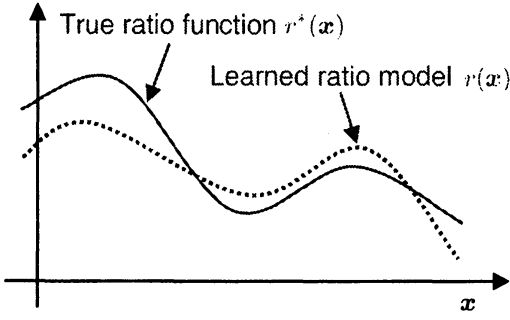
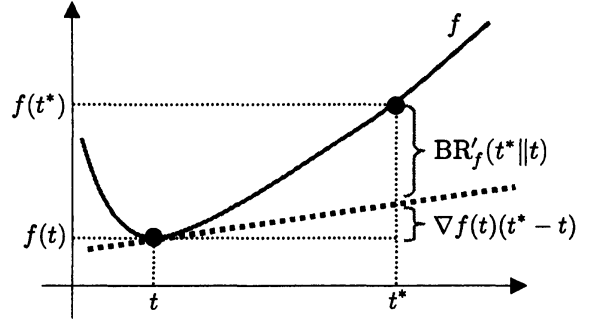


Figure 2: The idea of ratio matching.

Figure 3: Bregman divergence $BR'_f(t^*||t)$.

Here the discrepancy from the true density ratio r^* to a density ratio model r is measured using the BR divergence as

$$BR'_f(r^*||r) := \int p_{de}^*(\mathbf{x}) \left(f(r^*(\mathbf{x})) - f(r(\mathbf{x})) - \nabla f(r(\mathbf{x}))(r^*(\mathbf{x}) - r(\mathbf{x})) \right) d\mathbf{x}. \quad (9)$$

A motivation for this choice is that the BR divergence allows one to directly obtain an empirical approximation for any f . Indeed,

$$\begin{aligned} BR'_f(r^*||r) &= C - \int p_{de}^*(\mathbf{x}) f(r(\mathbf{x})) d\mathbf{x} - \int p_{de}^*(\mathbf{x}) \nabla f(r(\mathbf{x})) r^*(\mathbf{x}) d\mathbf{x} + \int p_{de}^*(\mathbf{x}) \nabla f(r(\mathbf{x})) r(\mathbf{x}) d\mathbf{x} \\ &= C + BR_f(r), \end{aligned}$$

where $C := \int p_{de}^*(\mathbf{x}) f(r^*(\mathbf{x})) d\mathbf{x}$ is a constant independent of r , and

$$BR_f(r) := \int p_{de}^*(\mathbf{x}) \nabla f(r(\mathbf{x})) r(\mathbf{x}) d\mathbf{x} - \int p_{de}^*(\mathbf{x}) f(r(\mathbf{x})) d\mathbf{x} - \int p_{nu}^*(\mathbf{x}) \nabla f(r(\mathbf{x})) d\mathbf{x}. \quad (10)$$

Thus, an empirical approximation $\widehat{BR}_f(r)$ of $BR_f(r)$ is given by

$$\widehat{BR}_f(r) := \frac{1}{n_{de}} \sum_{j=1}^{n_{de}} \nabla f(r(\mathbf{x}_j^{de})) r(\mathbf{x}_j^{de}) - \frac{1}{n_{de}} \sum_{j=1}^{n_{de}} f(r(\mathbf{x}_j^{de})) - \frac{1}{n_{nu}} \sum_{i=1}^{n_{nu}} \nabla f(r(\mathbf{x}_i^{nu})). \quad (11)$$

Below, ratio matching methods under the Kullback-Leibler divergence (Sugiyama et al., 2008) and the squared distance (Kanamori et al., 2009a) are explained.

4.2 Unnormalized Kullback-Leibler Divergence

In this section, a ratio matching method under the *unnormalized Kullback-Leibler (UKL) divergence* is explained.

4.2.1 Criterion

When $f(t) = t \log t - t$, the BR divergence is reduced to the UKL divergence:

$$UKL'(t^*||t) := t^* \log \frac{t^*}{t} - t^* + t.$$

Following Eqs.(10) and (11), let us denote UKL without an irrelevant constant term by $\text{UKL}(r)$ and its empirical approximation by $\widehat{\text{UKL}}(r)$:

$$\begin{aligned}\text{UKL}(r) &:= \int p_{\text{de}}^*(\mathbf{x})r(\mathbf{x})d\mathbf{x} - \int p_{\text{nu}}^*(\mathbf{x})\log r(\mathbf{x})d\mathbf{x}, \\ \widehat{\text{UKL}}(r) &:= \frac{1}{n_{\text{de}}} \sum_{j=1}^{n_{\text{de}}} r(\mathbf{x}_j^{\text{de}}) - \frac{1}{n_{\text{nu}}} \sum_{i=1}^{n_{\text{nu}}} \log r(\mathbf{x}_i^{\text{nu}}).\end{aligned}$$

The density ratio model r is learned so that $\widehat{\text{UKL}}(r)$ is minimized. Here, we further impose that the ratio model $r(\mathbf{x})$ is non-negative for all \mathbf{x} in \mathcal{X} and is normalized at $\{\mathbf{x}_j^{\text{de}}\}_{j=1}^{n_{\text{de}}}$:

$$\frac{1}{n_{\text{de}}} \sum_{j=1}^{n_{\text{de}}} r(\mathbf{x}_j^{\text{de}}) = 1.$$

Then the optimization criterion is reduced to as follows.

$$\max_r \frac{1}{n_{\text{nu}}} \sum_{i=1}^{n_{\text{nu}}} \log r(\mathbf{x}_i^{\text{nu}}) \quad \text{s.t.} \quad \frac{1}{n_{\text{de}}} \sum_{j=1}^{n_{\text{de}}} r(\mathbf{x}_j^{\text{de}}) = 1 \quad \text{and} \quad r(\mathbf{x}) \geq 0 \quad \text{for all } \mathbf{x} \in \mathcal{X}.$$

This is called the *KL importance estimation procedure* (KLIEP). Note that the same objective function can be obtained from an empirical approximation of the KL divergence from $p_{\text{nu}}^*(\mathbf{x})$ to $r(\mathbf{x})p_{\text{de}}^*(\mathbf{x})$ (Sugiyama et al., 2008).

Below, we describe how the KLIEP formulation can be implemented for linear and kernel models. Note that the KLIEP idea can be applied to various models such as *log-linear models* (Tsuboi et al., 2009), *Gaussian mixture models* (Yamada & Sugiyama, 2009), and *mixtures of probabilistic principal component analyzers* (Yamada et al., 2010).

4.2.2 Linear and Kernel Models

For the linear density-ratio model (2), the KLIEP optimization problem is expressed as follows (Sugiyama et al., 2008):

$$\max_{\boldsymbol{\theta} \in \mathbb{R}^b} \frac{1}{n_{\text{nu}}} \sum_{i=1}^{n_{\text{nu}}} \log(\boldsymbol{\psi}(\mathbf{x}_i^{\text{nu}})^\top \boldsymbol{\theta}) \quad \text{s.t.} \quad \bar{\boldsymbol{\psi}}_{\text{de}}^\top \boldsymbol{\theta} = 1 \quad \text{and} \quad \boldsymbol{\theta} \geq \mathbf{0}_b,$$

where $\bar{\boldsymbol{\psi}}_{\text{de}} := \frac{1}{n_{\text{de}}} \sum_{j=1}^{n_{\text{de}}} \boldsymbol{\psi}(\mathbf{x}_j^{\text{de}})$, and the inequality for vectors is applied in the element-wise manner. Since the above optimization problem is *convex* (i.e., the objective function to be maximized is concave and the feasible set is convex), there exists the unique global optimum solution. A pseudo code of KLIEP for linear models is described in Figure 4. As can be confirmed from the pseudo code, the denominator samples $\{\mathbf{x}_j^{\text{de}}\}_{j=1}^{n_{\text{de}}}$ appear only in terms of the basis-transformed mean $\bar{\boldsymbol{\psi}}_{\text{de}}$. Thus, KLIEP is computationally very efficient even when the number n_{de} of denominator samples is very large.

4.2.3 Basis Function Design

The performance of KLIEP depends on the choice of the basis functions $\boldsymbol{\psi}(\mathbf{x})$. As explained below, the use of the following Gaussian kernel model would be reasonable:

$$r(\mathbf{x}) = \sum_{\ell=1}^{n_{\text{nu}}} \theta_\ell K(\mathbf{x}, \mathbf{x}_\ell^{\text{nu}}), \tag{12}$$

```

Input: Data samples  $D^{\text{nu}} = \{\mathbf{x}_i^{\text{nu}}\}_{i=1}^{n_{\text{nu}}}$  and  $D^{\text{de}} = \{\mathbf{x}_j^{\text{de}}\}_{j=1}^{n_{\text{de}}}$ ,
        and basis functions  $\psi(\mathbf{x})$ 
Output: Density ratio estimator  $\hat{r}(\mathbf{x})$ 

 $\Psi_{\text{nu}} \leftarrow (\psi(\mathbf{x}_1^{\text{nu}}), \dots, \psi(\mathbf{x}_{n_{\text{nu}}}^{\text{nu}}))^{\top}$ ;
 $\bar{\psi}_{\text{de}} \leftarrow \frac{1}{n_{\text{de}}} \sum_{j=1}^{n_{\text{de}}} \psi(\mathbf{x}_j^{\text{de}})$ ;
Initialize  $\theta (> \mathbf{0}_b)$  and  $\varepsilon$  ( $0 < \varepsilon \ll 1$ );
Repeat until convergence
     $\theta \leftarrow \theta + \varepsilon \Psi_{\text{nu}}^{\top} (\mathbf{1}_{n_{\text{nu}}} / \Psi_{\text{nu}} \theta)$ ; % Gradient ascent
     $\theta \leftarrow \theta + (1 - \bar{\psi}_{\text{de}}^{\top} \theta) \bar{\psi}_{\text{de}} / (\bar{\psi}_{\text{de}}^{\top} \bar{\psi}_{\text{de}})$ ; % Constraint satisfaction
     $\theta \leftarrow \max(\mathbf{0}_b, \theta)$ ; % Constraint satisfaction
     $\theta \leftarrow \theta / (\bar{\psi}_{\text{de}}^{\top} \theta)$ ; % Constraint satisfaction
end
 $\hat{r}(\mathbf{x}) \leftarrow \psi(\mathbf{x})^{\top} \theta$ ;

```

Figure 4: Pseudo code of KLIEP. ‘./’ indicates the element-wise division and $^{\top}$ denotes the transpose. Inequalities and the ‘max’ operation for vectors are applied in the element-wise manner.

```

Input: Data samples  $D^{\text{nu}} = \{\mathbf{x}_i^{\text{nu}}\}_{i=1}^{n_{\text{nu}}}$  and  $D^{\text{de}} = \{\mathbf{x}_j^{\text{de}}\}_{j=1}^{n_{\text{de}}}$ ,
        and a set of basis function candidates  $\{\psi_m(\mathbf{x})\}_{m=1}^M$ 
Output: Density ratio estimator  $\hat{r}(\mathbf{x})$ 

Split  $D^{\text{nu}}$  into  $T$  disjoint subsets  $\{D_t^{\text{nu}}\}_{t=1}^T$ ;
for each model candidate  $m = 1, \dots, M$ 
    for each split  $t = 1, \dots, T$ 
         $\hat{r}_t(\mathbf{x}) \leftarrow \text{KLIEP}(D^{\text{nu}} \setminus D_t^{\text{nu}}, D^{\text{de}}, \psi(\mathbf{x}))$ ;
         $\widehat{\text{UKL}}_t(m) \leftarrow \frac{1}{|D_t^{\text{nu}}|} \sum_{\mathbf{x} \in D_t^{\text{nu}}} \log \hat{r}_t(\mathbf{x})$ ;
    end
     $\widehat{\text{UKL}}(m) \leftarrow \frac{1}{T} \sum_{t=1}^T \widehat{\text{UKL}}_t(m)$ ;
end
 $\hat{m} \leftarrow \text{argmax}_m \widehat{\text{UKL}}(m)$ ;
 $\hat{r}(\mathbf{x}) \leftarrow \text{KLIEP}(D^{\text{nu}}, D^{\text{de}}, \psi_{\hat{m}}(\mathbf{x}))$ ;

```

Figure 5: Pseudo code of CV for KLIEP.

where $K(\mathbf{x}, \mathbf{x}')$ is the Gaussian kernel (4). The reason why the numerator samples $\{\mathbf{x}_i^{\text{nu}}\}_{i=1}^{n_{\text{nu}}}$, not the denominator samples $\{\mathbf{x}_j^{\text{de}}\}_{j=1}^{n_{\text{de}}}$, are chosen as the Gaussian centers is as follows. By definition, the density ratio $r^*(\mathbf{x})$ tends to take large values if $p_{\text{de}}^*(\mathbf{x})$ is small and $p_{\text{nu}}^*(\mathbf{x})$ is large. Conversely, $r^*(\mathbf{x})$ tends to be small (i.e., close to zero) if $p_{\text{de}}^*(\mathbf{x})$ is large and $p_{\text{nu}}^*(\mathbf{x})$ is small. When a non-negative function is approximated by a Gaussian kernel model, many kernels may be needed in the region where the output of the target function is large. On the other hand, only a small number of kernels would be enough in the region where the output of the target function is close to zero (see Figure 6). Following this heuristic, many kernels are allocated in the region where $p_{\text{nu}}^*(\mathbf{x})$ takes large values, which can be achieved by setting the Gaussian centers at $\{\mathbf{x}_i^{\text{nu}}\}_{i=1}^{n_{\text{nu}}}$.

Alternatively, we may locate $(n_{\text{nu}} + n_{\text{de}})$ Gaussian kernels at both $\{\mathbf{x}_i^{\text{nu}}\}_{i=1}^{n_{\text{nu}}}$ and $\{\mathbf{x}_j^{\text{de}}\}_{j=1}^{n_{\text{de}}}$. However, this seems not to further improve the performance, but slightly increases the computa-

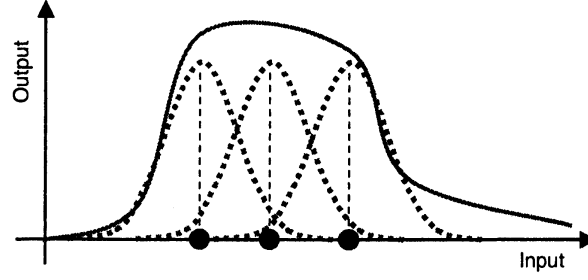


Figure 6: Heuristic of Gaussian kernel allocation.

tional cost. When n_{nu} is very large, just using all the numerator samples $\{\mathbf{x}_i^{\text{nu}}\}_{i=1}^{n_{\text{nu}}}$ as Gaussian centers is already computationally rather demanding. To ease this problem, a subset of $\{\mathbf{x}_i^{\text{nu}}\}_{i=1}^{n_{\text{nu}}}$ may be chosen in practice as Gaussian centers for computational efficiency, i.e.,

$$r(\mathbf{x}) = \sum_{\ell=1}^b \theta_{\ell} K(\mathbf{x}, \mathbf{c}_{\ell}),$$

where \mathbf{c}_{ℓ} is a template point randomly chosen from $\{\mathbf{x}_i^{\text{nu}}\}_{i=1}^{n_{\text{nu}}}$ and $b \in \{1, \dots, n_{\text{nu}}\}$ is a prefixed number.

4.2.4 Model Selection

Model selection of KLIEP is possible based on a variant of CV. More specifically, the numerator samples $D^{\text{nu}} = \{\mathbf{x}_i^{\text{nu}}\}_{i=1}^{n_{\text{nu}}}$ are divided into T disjoint subsets $\{D_t^{\text{nu}}\}_{t=1}^T$. Then a KLIEP solution $\hat{r}_t(\mathbf{x})$ is obtained using $D^{\text{nu}} \setminus D_t^{\text{nu}}$ (i.e., all numerator samples without D_t^{nu}) and D^{de} , and its UKL value for the hold-out samples D_t^{nu} is computed:

$$\widetilde{\text{UKL}}_t := \frac{1}{|D_t^{\text{nu}}|} \sum_{\mathbf{x}^{\text{nu}} \in D_t^{\text{nu}}} \log \hat{r}_t(\mathbf{x}^{\text{nu}}).$$

This procedure is repeated for $t = 1, \dots, T$, and the average of the above hold-out UKL values over all t is computed.

$$\widetilde{\text{UKL}} := \frac{1}{T} \sum_{t=1}^T \widetilde{\text{UKL}}_t.$$

Then the model that maximizes $\widetilde{\text{UKL}}$ is chosen.

A pseudo code of CV for KLIEP is summarized in Figure 5. A MATLAB[®] implementation of the entire KLIEP algorithm is available from

<http://sugiyama-www.cs.titech.ac.jp/~sugi/software/KLIEP/>

4.3 Squared Distance

Here, ratio matching methods under the *squared (SQ) distance* are described.

4.3.1 Criterion

When $f(t) = \frac{1}{2}(t-1)^2$, the BR divergence is reduced to the SQ distance:

$$\text{SQ}'(t^*||t) := \frac{1}{2}(t^* - t)^2.$$

Following Eqs.(10) and (11), let us denote SQ without an irrelevant constant term by $\text{SQ}(r)$ and its empirical approximation by $\widehat{\text{SQ}}(r)$:

$$\begin{aligned} \text{SQ}(r) &:= \frac{1}{2} \int p_{\text{de}}^*(\mathbf{x}) r(\mathbf{x})^2 d\mathbf{x} - \int p_{\text{nu}}^*(\mathbf{x}) r(\mathbf{x}) d\mathbf{x}, \\ \widehat{\text{SQ}}(r) &:= \frac{1}{2n_{\text{de}}} \sum_{j=1}^{n_{\text{de}}} r(\mathbf{x}_j^{\text{de}})^2 - \frac{1}{n_{\text{nu}}} \sum_{i=1}^{n_{\text{nu}}} r(\mathbf{x}_i^{\text{nu}}). \end{aligned}$$

Here, we focus on using the linear density-ratio model (2). Since this is the same model as KLIEP for linear models, the basis design heuristic introduced in Section 4.2.3 may also be used here. For the linear density-ratio model (2), $\widehat{\text{SQ}}$ is expressed as

$$\widehat{\text{SQ}}(\boldsymbol{\theta}) := \frac{1}{2} \boldsymbol{\theta}^\top \widehat{\mathbf{H}} \boldsymbol{\theta} - \widehat{\mathbf{h}}^\top \boldsymbol{\theta},$$

where

$$\widehat{\mathbf{H}} := \frac{1}{n_{\text{de}}} \sum_{j=1}^{n_{\text{de}}} \boldsymbol{\psi}(\mathbf{x}_j^{\text{de}}) \boldsymbol{\psi}(\mathbf{x}_j^{\text{de}})^\top \quad \text{and} \quad \widehat{\mathbf{h}} := \frac{1}{n_{\text{nu}}} \sum_{i=1}^{n_{\text{nu}}} \boldsymbol{\psi}(\mathbf{x}_i^{\text{nu}}).$$

4.3.2 Constrained Formulation

We impose non-negativity constraint $\boldsymbol{\theta} \geq \mathbf{0}_b$ when minimizing $\widehat{\text{SQ}}$. Then the optimization problem is expressed as follows.

$$\max_{\boldsymbol{\theta} \in \mathbb{R}^b} \frac{1}{2} \boldsymbol{\theta}^\top \widehat{\mathbf{H}} \boldsymbol{\theta} - \widehat{\mathbf{h}}^\top \boldsymbol{\theta} + \lambda \mathbf{1}_b^\top \boldsymbol{\theta} \quad \text{s.t.} \quad \boldsymbol{\theta} \geq \mathbf{0}_b, \quad (13)$$

where $\lambda (\geq 0)$ is the regularization parameter, and the constraint is imposed in order to guarantee the non-negativity of the density ratio estimator (given that the basis functions are non-negative). Together with the non-negativity constraint, the term $\mathbf{1}_b^\top \boldsymbol{\theta}$ works as the ℓ_1 -regularizer:

$$\mathbf{1}_b^\top \boldsymbol{\theta} = \|\boldsymbol{\theta}\|_1 := \sum_{\ell=1}^b |\theta_\ell|.$$

This formulation is called *least-squares importance fitting* (LSIF) (Kanamori et al., 2009a). The LSIF optimization problem is a convex quadratic program. Therefore, the unique global optimal solution can be computed by a standard optimization package.

We can also use the ℓ_2 -regularizer $\boldsymbol{\theta}^\top \boldsymbol{\theta}$, instead of the ℓ_1 -regularizer $\mathbf{1}_b^\top \boldsymbol{\theta}$, without changing the computational property. However, using the ℓ_1 -regularizer would be more advantageous since the solution tends to be *sparse* (Williams, 1995; Tibshirani, 1996; Chen et al., 1998). Furthermore, as explained in Section 4.3.3, the use of the ℓ_1 -regularizer allows one to compute the entire *regularization path* efficiently.

Model selection of LSIF (i.e., the choice of the basis functions and the regularization parameter) is possible by CV based on the SQ distance. More specifically, the numerator and

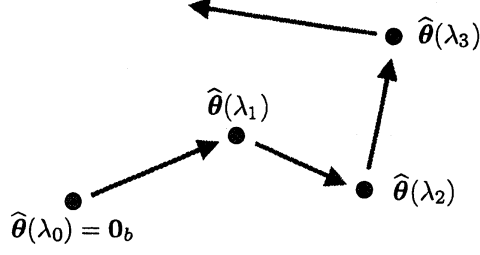


Figure 7: Regularization path tracking of LSIF. The solution $\hat{\theta}(\lambda)$ is shown to be piecewise-linear in the parameter space as a function of λ . Starting from $\lambda = \infty$, the trajectory of the solution is traced as λ is decreased to zero. When $\lambda \geq \lambda_0$ for some $\lambda_0 \geq 0$, the solution stays at the origin $\mathbf{0}_b$. When λ gets smaller than λ_0 , the solution departs from the origin. As λ is further decreased, for some λ_1 such that $0 \leq \lambda_1 \leq \lambda_0$, the solution goes straight to $\hat{\theta}(\lambda_1)$ with a constant ‘speed’. Then the solution path changes the direction and, for some λ_2 such that $0 \leq \lambda_2 \leq \lambda_1$, the solution is headed straight for $\hat{\theta}(\lambda_2)$ with a constant speed as λ is further decreased. This process is repeated until λ reaches zero.

denominator samples $D^{\text{nu}} = \{\mathbf{x}_i^{\text{nu}}\}_{i=1}^{n_{\text{nu}}}$ and $D^{\text{de}} = \{\mathbf{x}_j^{\text{de}}\}_{j=1}^{n_{\text{de}}}$ are divided into T disjoint subsets $\{D_t^{\text{nu}}\}_{t=1}^T$ and $\{D_t^{\text{de}}\}_{t=1}^T$, respectively. Then a density ratio estimator $\hat{r}_t(\mathbf{x})$ is obtained using $D^{\text{nu}} \setminus D_t^{\text{nu}}$ and $D^{\text{de}} \setminus D_t^{\text{de}}$ (i.e., all samples without D_t^{nu} and D_t^{de}), and its SQ value for the hold-out samples D_t^{nu} and D_t^{de} is computed:

$$\widetilde{\text{SQ}}_t := \frac{1}{2|D_t^{\text{nu}}|} \sum_{\mathbf{x}^{\text{nu}} \in D_t^{\text{nu}}} r(\mathbf{x}^{\text{nu}})^2 - \frac{1}{|D_t^{\text{de}}|} \sum_{\mathbf{x}^{\text{de}} \in D_t^{\text{de}}} r(\mathbf{x}^{\text{de}}).$$

This procedure is repeated for $t = 1, \dots, T$, and the average of the above hold-out SQ values is computed.

$$\widetilde{\text{SQ}} := \frac{1}{T} \sum_{t=1}^T \widetilde{\text{SQ}}_t.$$

Then the model that minimizes $\widetilde{\text{SQ}}$ is chosen.

For LSIF, an *information criterion* (Akaike, 1974) is also available for model selection (Kanamori et al., 2009a).

4.3.3 Entire Regularization Path

The LSIF solution $\hat{\theta}$ is shown to be *piecewise-linear* with respect to the regularization parameter λ (see Figure 7). Thus, the *regularization path* (i.e., solutions for all λ) can be computed efficiently based on the *parametric optimization technique* (Best, 1982; Efron et al., 2004; Hastie et al., 2004).

A basic idea of regularization path tracking is to check the violation of the *Karush-Kuhn-Tucker (KKT) conditions* (Boyd & Vandenberghe, 2004)—which are necessary and sufficient for optimality of convex programs—when the regularization parameter λ is changed. A pseudo code of the regularization path tracking algorithm for LSIF is described in Figure 8. Thanks to the regularization path algorithm, LSIF is computationally efficient in model selection scenarios, where solutions for various λ are computed.

The pseudo code implies that we no longer need a quadratic programming solver for obtaining the solution of LSIF—just computing matrix inverses is sufficient. Furthermore, the regularization path algorithm is computationally more efficient when the solution is sparse, that is, most

```

Input:  $\widehat{H}$  and  $\widehat{h}$ 
Output: entire regularization path  $\widehat{\theta}(\lambda)$  for  $\lambda \geq 0$ 

 $\tau \leftarrow 0$ ;  $k \leftarrow \operatorname{argmax}_i \{\widehat{h}_i \mid i = 1, \dots, b\}$ ;  $\lambda_\tau \leftarrow \widehat{h}_k$ ;
 $\widehat{\mathcal{A}} \leftarrow \{1, \dots, b\} \setminus \{k\}$ ;  $\widehat{\theta}(\lambda_\tau) \leftarrow \mathbf{0}_b$ ;
While  $\lambda_\tau > 0$ 
     $\widehat{E} \leftarrow O_{|\widehat{\mathcal{A}}| \times b}$ ;
    For  $i = 1, \dots, |\widehat{\mathcal{A}}|$ 
         $\widehat{E}_{i, \widehat{j}_i} \leftarrow 1$ ;  $\% \widehat{\mathcal{A}} = \{\widehat{j}_1, \dots, \widehat{j}_{|\widehat{\mathcal{A}}|} \mid \widehat{j}_1 < \dots < \widehat{j}_{|\widehat{\mathcal{A}}|}\}$ 
    end
     $\widehat{G} \leftarrow \begin{pmatrix} \widehat{H} & -\widehat{E}^\top \\ -\widehat{E} & O_{|\widehat{\mathcal{A}}| \times |\widehat{\mathcal{A}}|} \end{pmatrix}$ ;  $u \leftarrow \widehat{G}^{-1} \begin{pmatrix} \widehat{h} \\ \mathbf{0}_{|\widehat{\mathcal{A}}|} \end{pmatrix}$ ;  $v \leftarrow \widehat{G}^{-1} \begin{pmatrix} \mathbf{1}_b \\ \mathbf{0}_{|\widehat{\mathcal{A}}|} \end{pmatrix}$ ;
    If  $v \leq \mathbf{0}_{b+|\widehat{\mathcal{A}}|}$   $\% \text{ the final interval}$ 
         $\lambda_{\tau+1} \leftarrow 0$ ;  $\widehat{\theta}(\lambda_{\tau+1}) \leftarrow (u_1, \dots, u_b)^\top$ ;
    else  $\% \text{ an intermediate interval}$ 
         $k \leftarrow \operatorname{argmax}_i \{u_i/v_i \mid v_i > 0, i = 1, \dots, b + |\widehat{\mathcal{A}}|\}$ ;  $\lambda_{\tau+1} \leftarrow \max\{0, u_k/v_k\}$ ;
         $\widehat{\theta}(\lambda_{\tau+1}) \leftarrow (u_1, \dots, u_b)^\top - \lambda_{\tau+1}(v_1, \dots, v_b)^\top$ ;
        If  $1 \leq k \leq b$ 
             $\widehat{\mathcal{A}} \leftarrow \widehat{\mathcal{A}} \cup \{k\}$ ;
        else
             $\widehat{\mathcal{A}} \leftarrow \widehat{\mathcal{A}} \setminus \{\widehat{j}_{k-b}\}$ ;
        end
    end
     $\tau \leftarrow \tau + 1$ ;
end

 $\widehat{\theta}(\lambda) \leftarrow \begin{cases} \mathbf{0}_b & \text{if } \lambda \geq \lambda_0 \\ \frac{\lambda_{\tau+1}-\lambda}{\lambda_{\tau+1}-\lambda_\tau} \widehat{\theta}(\lambda_\tau) + \frac{\lambda-\lambda_\tau}{\lambda_{\tau+1}-\lambda_\tau} \widehat{\theta}(\lambda_{\tau+1}) & \text{if } \lambda_{\tau+1} \leq \lambda \leq \lambda_\tau \end{cases}$ 

```

Figure 8: Pseudo code for computing the entire regularization path of LSIF. When the computation of \widehat{G}^{-1} is numerically unstable, we may add small positive diagonals to \widehat{H} for stabilization purposes.

of the elements are zero since the number of change points tends to be small for such sparse solutions. However, the regularization path tracking algorithm was found to be numerically rather unstable (Kanamori et al., 2009a).

An R implementation of LSIF is available from

<http://www.math.cm.is.nagoya-u.ac.jp/~kanamori/software/LSIF/>

4.3.4 Unconstrained Formulation

The regularization path tracking algorithm for LSIF was shown to suffer from a numerical problem, and therefore is not practically reliable. Here, a practical alternative to LSIF is introduced, which gives an approximate solution to LSIF in a computationally efficient and reliable manner (Kanamori et al., 2009a).

The approximation idea introduced here is very simple: the non-negativity constraint of the

parameters in the optimization problem (13) is ignored. This results in the following unconstrained optimization problem.

$$\min_{\beta \in \mathbb{R}^b} \left[\frac{1}{2} \beta^\top \widehat{\mathbf{H}} \beta - \widehat{\mathbf{h}}^\top \beta + \frac{\lambda}{2} \beta^\top \beta \right]. \quad (14)$$

In the above, a quadratic regularization term $\beta^\top \beta / 2$ was included instead of the linear one $\mathbf{1}_b^\top \theta$ since the linear penalty term does not work as a regularizer without the non-negativity constraint. Eq.(14) is an unconstrained convex quadratic program, and the solution can be analytically computed as

$$\tilde{\beta} = (\widehat{\mathbf{H}} + \lambda \mathbf{I}_b)^{-1} \widehat{\mathbf{h}},$$

where \mathbf{I}_b is the b -dimensional identity matrix. Since the non-negativity constraint $\beta \geq \mathbf{0}_b$ was dropped, some of the learned parameters could be negative. To compensate for this approximation error, the solution is modified as

$$\hat{\beta} = \max(\mathbf{0}_b, \tilde{\beta}),$$

where the ‘max’ operation for a pair of vectors is applied in the element-wise manner. This is the solution of the approximation method called *unconstrained LSIF* (uLSIF) (Kanamori et al., 2009a). An advantage of uLSIF is that the solution can be computed just by solving a system of linear equations. Therefore, its computation is stable when λ is not too small.

Due to the ℓ_2 -regularizer, the solution tends to be close to $\mathbf{0}_b$ to some extent. Thus, the effect of ignoring the non-negativity constraint may not be so critical. See Kanamori et al. (2009a) for theoretical and experimental error analysis.

4.3.5 Analytic Expression of Leave-One-Out Score

A practically important advantage of uLSIF over LSIF is that the score of *leave-one-out CV* (LOOCV) can be computed analytically (Kanamori et al., 2009a)—thanks to this property, the computational complexity for performing LOOCV is the same order as just computing a single solution.

In the current setup, two sets of samples, $\{\mathbf{x}_i^{\text{nu}}\}_{i=1}^{n_{\text{nu}}}$ and $\{\mathbf{x}_j^{\text{de}}\}_{j=1}^{n_{\text{de}}}$, generally have different sample size. For simplicity, we assume that $n_{\text{de}} \leq n_{\text{nu}}$ and the i -th numerator sample \mathbf{x}_i^{nu} and the i -th denominator sample \mathbf{x}_i^{de} are held out at the same time; the numerator samples $\{\mathbf{x}_i^{\text{nu}}\}_{i=n_{\text{de}}+1}^{n_{\text{nu}}}$ are always used for density ratio estimation. Note that this assumption is only for the sake of simplicity; the order of numerator samples can be changed without sacrificing the computational advantages.

Let $\hat{r}^{(i)}(\mathbf{x})$ be a density ratio estimate obtained without the i -th numerator sample \mathbf{x}_i^{nu} and the i -th denominator sample \mathbf{x}_i^{de} . Then the LOOCV score is expressed as

$$\text{LOOCV} = \frac{1}{n_{\text{de}}} \sum_{i=1}^{n_{\text{de}}} \left[\frac{1}{2} (\hat{r}^{(i)}(\mathbf{x}_i^{\text{de}}))^2 - \hat{r}^{(i)}(\mathbf{x}_i^{\text{nu}}) \right].$$

A trick to efficiently compute the LOOCV score is to use the *Sherman-Woodbury-Morrison formula* (Golub & Loan, 1996) for computing matrix inverses: for an invertible square matrix \mathbf{A} and vectors \mathbf{u} and \mathbf{v} such that $\mathbf{v}^\top \mathbf{A}^{-1} \mathbf{u} \neq -1$,

$$(\mathbf{A} + \mathbf{u} \mathbf{v}^\top)^{-1} = \mathbf{A}^{-1} - \frac{\mathbf{A}^{-1} \mathbf{u} \mathbf{v}^\top \mathbf{A}^{-1}}{1 + \mathbf{v}^\top \mathbf{A}^{-1} \mathbf{u}}.$$

Input: $\{\mathbf{x}_i^{\text{nu}}\}_{i=1}^{n_{\text{nu}}}$ and $\{\mathbf{x}_j^{\text{de}}\}_{j=1}^{n_{\text{de}}}$
Output: $\hat{r}(\mathbf{x})$

$b \leftarrow \min(100, n_{\text{nu}}); \quad n \leftarrow \min(n_{\text{nu}}, n_{\text{de}});$
 Randomly choose b centers $\{\mathbf{c}_\ell\}_{\ell=1}^b$ from $\{\mathbf{x}_i^{\text{nu}}\}_{i=1}^{n_{\text{nu}}}$ without replacement;
For each candidate of Gaussian width σ

$$\hat{H}_{\ell, \ell'} \leftarrow \frac{1}{n_{\text{de}}} \sum_{j=1}^{n_{\text{de}}} \exp \left(-\frac{\|\mathbf{x}_j^{\text{de}} - \mathbf{c}_\ell\|^2 + \|\mathbf{x}_j^{\text{de}} - \mathbf{c}_{\ell'}\|^2}{2\sigma^2} \right) \text{ for } \ell, \ell' = 1, \dots, b;$$

$$\hat{h}_\ell \leftarrow \frac{1}{n_{\text{nu}}} \sum_{i=1}^{n_{\text{nu}}} \exp \left(-\frac{\|\mathbf{x}_i^{\text{nu}} - \mathbf{c}_\ell\|^2}{2\sigma^2} \right) \text{ for } \ell = 1, \dots, b;$$

$$X_{\ell, i}^{\text{nu}} \leftarrow \exp \left(-\frac{\|\mathbf{x}_i^{\text{nu}} - \mathbf{c}_\ell\|^2}{2\sigma^2} \right) \text{ for } i = 1, \dots, n \text{ and } \ell = 1, \dots, b;$$

$$X_{\ell, i}^{\text{de}} \leftarrow \exp \left(-\frac{\|\mathbf{x}_i^{\text{de}} - \mathbf{c}_\ell\|^2}{2\sigma^2} \right) \text{ for } i = 1, \dots, n \text{ and } \ell = 1, \dots, b;$$

For each candidate of regularization parameter λ

$$\hat{B} \leftarrow \hat{H} + \frac{\lambda(n_{\text{de}} - 1)}{n_{\text{de}}} \mathbf{I}_b;$$

$$B_0 \leftarrow \hat{B}^{-1} \hat{\mathbf{h}} \mathbf{1}_n^\top + \hat{B}^{-1} X^{\text{de}} \text{diag} \left(\frac{\hat{\mathbf{h}}^\top \hat{B}^{-1} X^{\text{de}}}{n_{\text{de}} \mathbf{1}_n^\top - \mathbf{1}_b^\top (X^{\text{de}} * \hat{B}^{-1} X^{\text{de}})} \right);$$

$$B_1 \leftarrow \hat{B}^{-1} X^{\text{nu}} + \hat{B}^{-1} X^{\text{de}} \text{diag} \left(\frac{\mathbf{1}_b^\top (X^{\text{nu}} * \hat{B}^{-1} X^{\text{de}})}{n_{\text{de}} \mathbf{1}_n^\top - \mathbf{1}_b^\top (X^{\text{de}} * \hat{B}^{-1} X^{\text{de}})} \right);$$

$$B_2 \leftarrow \max \left(\mathbf{O}_{b \times n}, \frac{n_{\text{de}} - 1}{n_{\text{de}}(n_{\text{nu}} - 1)} (n_{\text{nu}} B_0 - B_1) \right);$$

$$\mathbf{w}_{\text{de}} \leftarrow (\mathbf{1}_b^\top (X^{\text{de}} * B_2))^\top; \quad \mathbf{w}_{\text{nu}} \leftarrow (\mathbf{1}_b^\top (X^{\text{nu}} * B_2))^\top;$$

$$\text{LOOCV}(\sigma, \lambda) \leftarrow \frac{\mathbf{w}_{\text{de}}^\top \mathbf{w}_{\text{de}}}{2n} - \frac{\mathbf{1}_n^\top \mathbf{w}_{\text{nu}}}{n};$$

end

end

$$(\hat{\sigma}, \hat{\lambda}) \leftarrow \underset{(\sigma, \lambda)}{\text{argmin}} \text{LOOCV}(\sigma, \lambda);$$

$$\tilde{H}_{\ell, \ell'} \leftarrow \frac{1}{n_{\text{de}}} \sum_{j=1}^{n_{\text{de}}} \exp \left(-\frac{\|\mathbf{x}_j^{\text{de}} - \mathbf{c}_\ell\|^2 + \|\mathbf{x}_j^{\text{de}} - \mathbf{c}_{\ell'}\|^2}{2\hat{\sigma}^2} \right) \text{ for } \ell, \ell' = 1, \dots, b;$$

$$\tilde{h}_\ell \leftarrow \frac{1}{n_{\text{nu}}} \sum_{i=1}^{n_{\text{nu}}} \exp \left(-\frac{\|\mathbf{x}_i^{\text{nu}} - \mathbf{c}_\ell\|^2}{2\hat{\sigma}^2} \right) \text{ for } \ell = 1, \dots, b;$$

$$\hat{\alpha} \leftarrow \max(\mathbf{0}_b, (\tilde{H} + \hat{\lambda} \mathbf{I}_b)^{-1} \tilde{\mathbf{h}});$$

$$\hat{w}(\mathbf{x}) \leftarrow \sum_{\ell=1}^b \hat{\alpha}_\ell \exp \left(-\frac{\|\mathbf{x} - \mathbf{c}_\ell\|^2}{2\hat{\sigma}^2} \right);$$

Figure 9: Pseudo code of uLSIF algorithm with LOOCV. $B * B'$ denotes the element-wise multiplication of matrices B and B' of the same size, that is, the (i, j) -th element is given by $B_{i,j} B'_{i,j}$. For n -dimensional vectors \mathbf{b} and \mathbf{b}' , $\text{diag}(\frac{\mathbf{b}}{\mathbf{b}'})$ denotes the $n \times n$ diagonal matrix with i -th diagonal element b_i/b'_i .

A pseudo code of uLSIF with LOOCV-based model selection is summarized in Figure 9. Note that the basis design heuristic explained in Section 4.2.3 is used in the pseudo code, but the analytic form of the LOOCV score is available for any basis functions.

A MATLAB[®] implementation of uLSIF is available from

<http://sugiyama-www.cs.titech.ac.jp/~sugi/software/uLSIF/>

and an R implementation of uLSIF is available from

<http://www.math.cm.is.nagoya-u.ac.jp/~kanamori/software/LSIF/>

4.4 Remarks

Density ratio estimation by ratio matching can successfully avoid density estimation. Furthermore, CV based on the target divergence functional is available for model selection.

We have described the ratio matching methods for the UKL divergence and the SQ distance. The UKL method (KLIEP) is applicable to a variety of models such as linear/kernel models (Sugiyama et al., 2008), log-linear models (Tsuboi et al., 2009), mixtures of Gaussians (Yamada & Sugiyama, 2009), and mixtures of probabilistic principal component analyzers (Yamada et al., 2010). On the other hand, the SQ methods are computationally more efficient. The constrained method (LSIF) for the ℓ_1 -regularizer is equipped with a regularization path tracking algorithm. Furthermore, its unconstrained variant (uLSIF) allows one to compute the density ratio estimator *analytically*; the leave-one-out CV score can also be computed in a closed form. Thus, the overall computation of uLSIF including model selection is highly efficient (Kanamori et al., 2009a).

The fact that uLSIF has an analytic-form solution is actually very useful beyond its computational efficiency. When one wants to optimize some criterion defined using a density ratio estimate (e.g., *mutual information*, Cover & Thomas, 2006), the analytic-form solution of uLSIF allows one to compute the *derivative* of the target criterion analytically. Then one can develop, e.g., gradient-based algorithms and (quasi-) Newton algorithms for optimization. This property can be successfully utilized, e.g., in identifying the central subspace in *sufficient dimensionality reduction* (Suzuki & Sugiyama, 2010), finding independent components in *independent component analysis* (Suzuki & Sugiyama, 2009), performing dependence minimizing regression in *causal inference* (Yamada & Sugiyama, 2010), and identifying the hetero-distributional subspace in *direct density ratio estimation with dimensionality reduction* (Sugiyama et al., 2010a).

The ratio matching approach can also be characterized as divergence estimation. Let f be a convex function such that $f(1) = 0$. Then the *Ali-Silver-Csiszár (ASC) divergence* associated with f from p_{de}^* to p_{nu}^* is defined as follows (Ali & Silvey, 1966; Csiszár, 1967):

$$\text{ASC}_f(p_{\text{de}}^* \| p_{\text{nu}}^*) := \int p_{\text{de}}^*(\mathbf{x}) f\left(\frac{p_{\text{nu}}^*(\mathbf{x})}{p_{\text{de}}^*(\mathbf{x})}\right) d\mathbf{x}.$$

Let f^* be the *Legendre-Fenchel dual* of f (Rockafellar, 1970):

$$f^*(s) := \sup_{t^*} [t^* s - f(t^*)].$$

The convexity of f implies

$$f(t^*) \geq f(t) + (t^* - t) \nabla f(t),$$

where the equality holds if and only if $t^* = t$ (see Figure 3 again). Thus, for $s = \nabla f(t)$, $f^*(s)$ is expressed as

$$\begin{aligned} f^*(\nabla f(t)) &= \sup_{t^*} [t^* \nabla f(t) - f(t^*)] \\ &= t \nabla f(t) - f(t). \end{aligned}$$

Then, the BR divergence associated with f from r^* to r without an irrelevant constant (see Eq.(10)) can be expressed in terms of f^* as

$$\text{BR}_f(r^* \| r) = \int p_{\text{de}}^*(\mathbf{x}) \left(f^*(\nabla f(r(\mathbf{x}))) - \nabla f(r(\mathbf{x})) r^*(\mathbf{x}) \right) d\mathbf{x}. \quad (15)$$

Eq.(15) is minimized with respect to r if and only if $r = r^*$ (Nguyen et al., 2010):

$$\min_r \text{BR}_f(r^* \| r) = \int p_{\text{de}}^*(\mathbf{x}) f(r^*(\mathbf{x})) d\mathbf{x}.$$

Consequently, the ASC divergence can be approximated as

$$\begin{aligned} \text{ASC}_f(p_{\text{de}}^* \| p_{\text{nu}}^*) &= \min_r \text{BR}_f(r^* \| r) \\ &\approx \min_r \left[\int \frac{1}{n_{\text{de}}} \sum_{j=1}^{n_{\text{de}}} f^*(\nabla f(r(\mathbf{x}_j^{\text{de}}))) - \frac{1}{n_{\text{nu}}} \sum_{i=1}^{n_{\text{nu}}} \nabla f(r(\mathbf{x}_i^{\text{nu}})) \right]. \end{aligned}$$

This agrees with the ASC-divergence estimator proposed in Nguyen et al. (2010).

5 Conclusions

In this paper, we provided a comprehensive review of density ratio estimation methods, including the moment matching approach (Section 2), the probabilistic classification approach (Section 3), and the ratio matching approach (Section 4). Through extensive experiments, these methods were shown to outperform the naive approach of taking the ratio of kernel density estimators (Sugiyama et al., 2008; Gretton et al., 2009; Kanamori et al., 2009a; Hido et al., 2010).

Theoretical analysis of these direct density-ratio estimators has also been carried out. For example, in Kanamori et al. (2010), the accuracy of (A) the ratio of maximum likelihood density estimators, (B) probabilistic classification with logistic regression, and (C) ratio matching under the Kullback-Leibler divergence has been theoretically compared in the parametric setup. The paper showed that, when the numerator and denominator densities are known to be members of the exponential family, (A) is better than (B) and (B) is better than (C). On the other hand, once the model assumption is violated, (C) was shown to be better than (A) and (B). Thus, in practical situations where no exact model is available, (C) would be the most promising approach to density ratio estimation.

For non-parametric cases, the convergence rate of the infinite-order moment matching approach (Gretton et al., 2009), ratio matching under the Kullback-Leibler divergence (Sugiyama et al., 2008; Nguyen et al., 2010), and ratio matching under the squared distance (Kanamori et al., 2009b) has been elucidated. However, it seems to be an open research topic to theoretically prove that these direct density-ratio estimators are really superior to the naive approach of taking the ratio of non-parametric density estimators.

Finally, the performance of density ratio estimation in high-dimensional problems can be further improved by *dimensionality reduction*. More specifically, density ratio estimation is

carried out only in a subspace in which the numerator and denominator densities are significantly different. Such approaches have been explored recently (Sugiyama et al., 2010b; Sugiyama et al., 2010a), and would be a promising direction for further improving the estimation accuracy of density ratios.

Acknowledgment

MS was supported by SCAT, AOARD, and the JST PRESTO program. TS was supported by Global COE Program “The research and training center for new development in mathematics”, MEXT, Japan. TK was supported by MEXT Grant-in-Aid for Young Scientists 20700251.

References

- Akaike, H. (1974). A new look at the statistical model identification. *IEEE Transactions on Automatic Control*, AC-19, 716–723.
- Ali, S. M., & Silvey, S. D. (1966). A general class of coefficients of divergence of one distribution from another. *Journal of the Royal Statistical Society, Series B*, 28, 131–142.
- Best, M. J. (1982). *An algorithm for the solution of the parametric quadratic programming problem* (Technical Report 82-24). Faculty of Mathematics, University of Waterloo.
- Bickel, S., Bogojeska, J., Lengauer, T., & Scheffer, T. (2008). Multi-task learning for HIV therapy screening. *Proceedings of 25th Annual International Conference on Machine Learning (ICML2008)* (pp. 56–63). Helsinki, Finland: Omnipress.
- Bickel, S., Brückner, M., & Scheffer, T. (2007). Discriminative learning for differing training and test distributions. *Proceedings of the 24th International Conference on Machine Learning* (pp. 81–88).
- Boyd, S., & Vandenberghe, L. (2004). *Convex optimization*. Cambridge, UK: Cambridge University Press.
- Bregman, L. M. (1967). The relaxation method of finding the common point of convex sets and its application to the solution of problems in convex programming. *USSR Computational Mathematics and Mathematical Physics*, 7, 200–217.
- Chen, S. S., Donoho, D. L., & Saunders, M. A. (1998). Atomic decomposition by basis pursuit. *SIAM Journal on Scientific Computing*, 20, 33–61.
- Cheng, K. F., & Chu, C. K. (2004). Semiparametric density estimation under a two-sample density ratio model. *Bernoulli*, 10, 583–604.
- Cover, T. M., & Thomas, J. A. (2006). *Elements of information theory*. Hoboken, NJ, USA: John Wiley & Sons, Inc. 2nd edition.
- Csiszár, I. (1967). Information-type measures of difference of probability distributions and indirect observation. *Studia Scientiarum Mathematicarum Hungarica*, 2, 229–318.
- Efron, B., Hastie, T., Tibshirani, R., & Johnstone, I. (2004). Least angle regression. *The Annals of Statistics*, 32, 407–499.
- Golub, G. H., & Loan, C. F. V. (1996). *Matrix computations*. Baltimore, MD: Johns Hopkins University Press.
- Gretton, A., Smola, A., Huang, J., Schmittfull, M., Borgwardt, K., & Schölkopf, B. (2009). Covariate shift by kernel mean matching. In J. Quiñero-Candela, M. Sugiyama, A. Schwaighofer and N. Lawrence (Eds.), *Dataset shift in machine learning*, chapter 8, 131–160. Cambridge, MA: MIT Press.
- Hastie, T., Rosset, S., Tibshirani, R., & Zhu, J. (2004). The entire regularization path for the support vector machine. *Journal of Machine Learning Research*, 5, 1391–1415.

- Hastie, T., Tibshirani, R., & Friedman, J. (2001). *The elements of statistical learning: Data mining, inference, and prediction*. New York: Springer.
- Hido, S., Tsuboi, Y., Kashima, H., Sugiyama, M., & Kanamori, T. (2008). Inlier-based outlier detection via direct density ratio estimation. *Proceedings of IEEE International Conference on Data Mining (ICDM2008)* (pp. 223–232). Pisa, Italy.
- Hido, S., Tsuboi, Y., Kashima, H., Sugiyama, M., & Kanamori, T. (2010). Statistical outlier detection using direct density ratio estimation. *Knowledge and Information Systems*. to appear.
- Kanamori, T., Hido, S., & Sugiyama, M. (2009a). A least-squares approach to direct importance estimation. *Journal of Machine Learning Research*, 10, 1391–1445.
- Kanamori, T., Suzuki, T., & Sugiyama, M. (2009b). *Condition number analysis of kernel-based density ratio estimation* (Technical Report). arXiv. <http://www.citebase.org/abstract?id=oai:arXiv.org:0912.2800>
- Kanamori, T., Suzuki, T., & Sugiyama, M. (2010). Theoretical analysis of density ratio estimation. *IEICE Transactions on Fundamentals of Electronics, Communications and Computer Sciences*. to appear.
- Kawahara, Y., & Sugiyama, M. (2009). Change-point detection in time-series data by direct density-ratio estimation. *Proceedings of 2009 SIAM International Conference on Data Mining (SDM2009)* (pp. 389–400). Sparks, Nevada, USA.
- Minka, T. P. (2007). *A comparison of numerical optimizers for logistic regression* (Technical Report). Microsoft Research. <http://research.microsoft.com/~minka/papers/logreg/minka-logreg.pdf>
- Nguyen, X., Wainwright, M. J., & Jordan, M. I. (2010). Estimating divergence functionals and the likelihood ratio by convex risk minimization. *IEEE Transactions on Information Theory*. to appear.
- Qin, J. (1998). Inferences for case-control and semiparametric two-sample density ratio models. *Biometrika*, 85, 619–639.
- Quiñonero-Candela, J., Sugiyama, M., Schwaighofer, A., & Lawrence, N. (Eds.). (2009). *Dataset shift in machine learning*. Cambridge, MA: MIT Press.
- Rockafellar, R. T. (1970). *Convex analysis*. Princeton, NJ, USA: Princeton University Press.
- Schölkopf, B., & Smola, A. J. (2002). *Learning with kernels*. Cambridge, MA: MIT Press.
- Shimodaira, H. (2000). Improving predictive inference under covariate shift by weighting the log-likelihood function. *Journal of Statistical Planning and Inference*, 90, 227–244.
- Smola, A., Song, L., & Teo, C. H. (2009). Relative novelty detection. *Twelfth International Conference on Artificial Intelligence and Statistics* (pp. 536–543).
- Steinwart, I. (2001). On the influence of the kernel on the consistency of support vector machines. *Journal of Machine Learning Research*, 2, 67–93.
- Sugiyama, M. (2010). Superfast-trainable multi-class probabilistic classifier by least-squares posterior fitting. *IEICE Transactions on Information and Systems*. submitted.
- Sugiyama, M., Hara, S., von Büna, P., Suzuki, T., Kanamori, T., & Kawanabe, M. (2010a). Direct density ratio estimation with dimensionality reduction. *Proceedings of 2010 SIAM International Conference on Data Mining (SDM2010)*. Columbus, Ohio, USA.
- Sugiyama, M., Kanamori, T., Suzuki, T., Hido, S., Sese, J., Takeuchi, I., & Wang, L. (2009). A density-ratio framework for statistical data processing. *IPSJ Transactions on Computer Vision and Applications*, 1, 183–208.
- Sugiyama, M., Kawanabe, M., & Chui, P. L. (2010b). Dimensionality reduction for density ratio estimation in high-dimensional spaces. *Neural Networks*, 23, 44–59.
- Sugiyama, M., Krauledat, M., & Müller, K.-R. (2007). Covariate shift adaptation by importance weighted cross validation. *Journal of Machine Learning Research*, 8, 985–1005.

- Sugiyama, M., & Müller, K.-R. (2005). Input-dependent estimation of generalization error under covariate shift. *Statistics & Decisions*, 23, 249–279.
- Sugiyama, M., Suzuki, T., & Kanamori, T. (2011). *Density ratio estimation in machine learning: A versatile tool for statistical data processing*. Cambridge, UK: Cambridge University Press. to appear.
- Sugiyama, M., Suzuki, T., Nakajima, S., Kashima, H., von Büna, P., & Kawanabe, M. (2008). Direct importance estimation for covariate shift adaptation. *Annals of the Institute of Statistical Mathematics*, 60, 699–746.
- Sugiyama, M., Takeuchi, I., Suzuki, T., Kanamori, T., Hachiya, H., & Okanohara, D. (2010c). Least-squares conditional density estimation. *IEICE Transactions on Information and Systems*, E93-D, 583–594.
- Sugiyama, M., von Büna, P., Kawanabe, M., & Müller, K.-R. (2010d). *Covariate shift adaptation: Towards machine learning in non-stationary environment*. Cambridge, MA: MIT Press. to appear.
- Suzuki, T., & Sugiyama, M. (2009). Estimating squared-loss mutual information for independent component analysis. *Independent Component Analysis and Signal Separation* (pp. 130–137). Berlin: Springer.
- Suzuki, T., & Sugiyama, M. (2010). Sufficient dimension reduction via squared-loss mutual information estimation. *Proceedings of The Thirteenth International Conference on Artificial Intelligence and Statistics (AISTATS2010)* (pp. 804–811). Chia Laguna, Sardinia, Italy.
- Suzuki, T., Sugiyama, M., Kanamori, T., & Sese, J. (2009a). Mutual information estimation reveals global associations between stimuli and biological processes. *BMC Bioinformatics*, 10, S52.
- Suzuki, T., Sugiyama, M., Sese, J., & Kanamori, T. (2008). Approximating mutual information by maximum likelihood density ratio estimation. *JMLR Workshop and Conference Proceedings* (pp. 5–20).
- Suzuki, T., Sugiyama, M., & Tanaka, T. (2009b). Mutual information approximation via maximum likelihood estimation of density ratio. *Proceedings of 2009 IEEE International Symposium on Information Theory (ISIT2009)* (pp. 463–467). Seoul, Korea.
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society, Series B*, 58, 267–288.
- Tsuboi, Y., Kashima, H., Hido, S., Bickel, S., & Sugiyama, M. (2009). Direct density ratio estimation for large-scale covariate shift adaptation. *Journal of Information Processing*, 17, 138–155.
- Williams, P. M. (1995). Bayesian regularization and pruning using a Laplace prior. *Neural Computation*, 7, 117–143.
- Yamada, M., & Sugiyama, M. (2009). Direct importance estimation with Gaussian mixture models. *IEICE Transactions on Information and Systems*, E92-D, 2159–2162.
- Yamada, M., & Sugiyama, M. (2010). Dependence minimizing regression with model selection for non-linear causal inference under non-Gaussian noise. *Proceedings of the Twenty-Third AAAI Conference on Artificial Intelligence (AAAI2010)*. to appear.
- Yamada, M., Sugiyama, M., Wichern, G., & Simin, J. (2010). Direct importance estimation with a mixture of probabilistic principal component analyzers. *IEICE Transactions on Information and Systems*. submitted.
- Zadrozny, B. (2004). Learning and evaluating classifiers under sample selection bias. *Proceedings of the Twenty-First International Conference on Machine Learning* (pp. 903–910). New York, NY: ACM Press.